

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-255172

(43)Date of publication of application : 01.10.1996

(51)Int.Cl.

G06F 17/30

(21)Application number : 07-083458

(71)Applicant : TOSHIBA CORP

(22)Date of filing : 16.03.1995

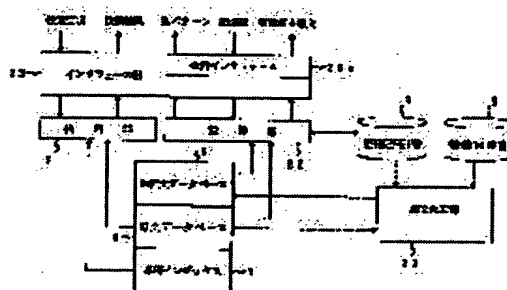
(72)Inventor : SAKAI TETSUYA
MIKE SEIJI
SUMITA KAZUO

(54) DOCUMENT RETRIEVAL SYSTEM

(57)Abstract:

PURPOSE: To lighten the burden of document retrieving operation on a user by extracting only desired extracted sentences or information from the whole text of a retrieved document and displaying it.

CONSTITUTION: The document retrieval system is equipped with an original text processing part 20 which sets plural kinds of sentence types for discriminating the contents of documents like an opinion and a proposal, generates extracted sentence data in sentence units classified by the types of the respective sentences from an original text database 6 storing original text data constituting the document, and stores the data as an extracted sentence database 5. The original text processing part 20 is equipped with a style decision part which extracts extracted sentence data corresponding to a specified sentence style and a shaping part which shapes the extracted sentence data into a specific style having, for example, conjunctions removed.



LEGAL STATUS

[Date of request for examination] 16.04.1999

[Date of sending the examiner's decision of rejection] 28.03.2000

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

BEST AVAILABLE COPY

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平8-255172

(43)公開日 平成8年(1996)10月1日

(51)Int.Cl. ⁶	識別記号	片内整理番号	F I	技術表示箇所
G 0 6 F 17/30		9194-5L	G 0 6 F 15/401	3 2 0 A
		9194-5L	15/403	3 8 0 D

審査請求 未請求 請求項の数6 F D (全 23 頁)

(21)出願番号 特願平7-83458

(22)出願日 平成7年(1995)3月16日

(71)出願人 000003078

株式会社東芝

神奈川県川崎市幸区堀川町72番地

(72)発明者 酒井 哲也

神奈川県川崎市幸区小向東芝町1番地 株
式会社東芝研究開発センター内

(72)発明者 三池 誠司

神奈川県川崎市幸区小向東芝町1番地 株
式会社東芝研究開発センター内

(72)発明者 住田 一男

神奈川県川崎市幸区小向東芝町1番地 株
式会社東芝研究開発センター内

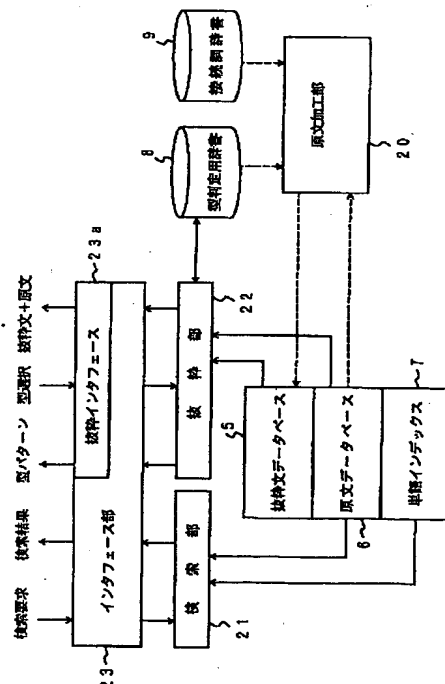
(74)代理人 弁理士 鈴江 武彦

(54)【発明の名称】 文書検索システム

(57)【要約】 (修正有)

【目的】 検索された文書の全文から所望の抜粋文や情報のみを抽出して表示できるようにして、利用者の文書検索作業に要する負荷の軽減化を図ることにある。

【構成】 文書を構成する原文データを格納した原文データベース6から、例えば意見、提言等のように文章の内容を識別するための複数種類の文の型を設定し、この各文の型に分類した文単位の抜粋文データを作成し、抜粋文データベース5として格納する原文加工部20を備えた文書検索システムである。原文加工部20は、指定された文の型に対応する抜粋文データを抽出する型判定部と、抜粋文データを例えば接続詞を除去したような所定の形式に整形する整形部とを備える。



1

【 特許請求の範囲】

【 請求項1 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、
前記原文データから予め設定した複数種の文の型に分類した文単位の抜粋文データを格納した抜粋文データベース手段と、
前記複数種の文の型から所望の文の型を選択する選択手段と、
この選択手段により選択された文の型に対応する抜粋文データを前記抜粋文データベース手段から抽出する抽出手段と、
この抽出手段により抽出された全ての抜粋文データを一覧的に表示する表示手段とを具備したことを特徴とする文書検索システム。

【 請求項2 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、
前記原文データから予め設定した複数種の文の型に分類した文単位の抜粋文データを作成する作成手段と、
この作成手段により作成された各抜粋文データを所定の形式に整形する整形手段と、
この整形手段により整形された各抜粋文データを格納した抜粋文データベース手段と、
前記複数種の文の型から所望の文の型を選択する選択手段と、
この選択手段により選択された文の型に対応する抜粋文データを前記抜粋文データベース手段から抽出する抽出手段と、
この抽出手段により抽出された全ての抜粋文データを一覧的に表示し、かつ抽出された前記各抜粋文データに対応する前記原文データを強調的に表示する表示手段とを具備したことを特徴とする文書検索システム。

【 請求項3 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、
前記原文データから予め設定した複数種の文の型に分類した文単位の抜粋文データを作成する作成手段と、
この作成手段により作成された各抜粋文データを所定の形式に整形する整形手段と、
この整形手段により整形された各抜粋文データを格納した抜粋文データベース手段と、
前記複数種の文の型から所望の文の型または全文表示モードを選択する選択手段と、
この選択手段により選択された文の型に対応する抜粋文データを前記抜粋文データベース手段から抽出する抽出手段と、

2

前記選択手段により 前記全文表示モードを選択された場合には、前記原文データベース手段から検索対象の文書の全文に対応する原文データを取出して出力する全文出力手段と、
前記抽出手段により抽出された全ての抜粋文データまたは前記全文出力手段により出力された原文データを表示する表示手段とを具備したことを特徴とする文書検索システム。

【 請求項4 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、
前記原文データから予め設定した複数種の文の型に分類した文単位の抜粋文データを作成する作成手段と、
この作成手段により作成された各抜粋文データを所定の形式に整形する整形手段と、
この整形手段により整形された各抜粋文データを格納した抜粋文データベース手段と、
前記複数種の文の型から所望の文の型を選択する選択手段と、
この選択手段により選択された文の型に対応する抜粋文データを前記抜粋文データベース手段から抽出する抽出手段と、
予め指定したキーワードに基づいて、前記抽出手段が抽出された抜粋文データの中で前記キーワードを含む特定の抜粋文データのみを選択して表示する表示手段とを具備したことを特徴とする文書検索システム。

【 請求項5 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、
前記原文データに対して文単位の言語解析処理を実行し、この言語解析処理により得られた言語情報を生成する言語情報生成手段と、
前記原文データから予め設定した複数種の文の型に分類した文単位の抜粋文データであって、前記言語情報を付加した抜粋文データを格納した抜粋文データベース手段と、
前記複数種の文の型から所望の文の型を選択する選択手段と、
この選択手段により選択された文の型に対応する抜粋文データを前記抜粋文データベース手段から抽出する抽出手段と、
この抽出手段により抽出された全ての抜粋文データを前記言語情報を伴って一覧的に表示する表示手段とを具備したことを特徴とする文書検索システム。

【 請求項6 】 検索要求に応じた文書を検索する文書検索システムにおいて、
検索対象の文書を構成する原文データを格納した原文データベース手段と、

10

20

30

40

50

3

前記原文データに対して言語解析処理を実行し、この言語解析処理により得られた言語情報に基づいた文の所定の構成要素を前記原文データから抽出する言語解析処理手段と、

この言語解析処理手段により抽出された前記所定の構成要素毎に予め設定した複数種の情報の型に分類した抜粋情報データを格納した抜粋情報データベース手段と、前記複数種の情報の型から所望の情報の型を選択する選択手段と、

この選択手段により選択された情報の型に対応する抜粋情報データを前記抜粋情報データベース手段から抽出する抽出手段と、

この抽出手段により抽出された全ての抜粋情報データを一覧的に表示する表示手段とを具備したことを特徴とする文書検索システム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、特に文書単位に文書情報を検索するシステムにおいて、文書の原文を加工処理して得られた抜粋文または抜粋情報を検索する機能を備えた文書検索システムに関する。

【0002】

【従来の技術】従来、フルテキストサーチ方式の文書検索システムは、文書データベースから必要な文書や情報を検索する場合に、文書を構成する文字または単語の全てを検索対象とするため、例えばキーワード検索方式と比較して、キーワードを付与する作業を必要せず、また利用者がそのキーワードに精通している必要がないなどの利点を備えている。

【0003】しかしながら、フルテキストサーチ方式では入力した単語や文章等の検索要求に対して、その検索要求に合致する文書数と文書データ量は膨大である。この方式では、検索処理により検索要求に合致する各文書の全文が、そのままの状態では表示画面上に表示される。利用者は、表示画面上で検索された各文書の全文を確認し、この中から必要な文書または情報(文章や単語からなる)を抽出することになる。

【0004】したがって、利用者は、検索された膨大な量の文書の原文に目を通す必要があるため、必要な文書や情報を得るまでに多大な時間と労力が必要となる。このような問題点を解消するためには、検索された文書の原文から必要な文や情報のみを抽出して、見易い形式で表示画面上に表示するシステムが望ましい。

【0005】近年、文書から自動的に抄録や要約を作成する技術が開発されており、この技術を利用して検索された文書中から重要文のみを抽出して表示するシステムが実現されている。また、文書中から例えば結論に相当する章または項の文章を抽出して表示するシステムも開発されている。

【0006】このようなシステムにおいても、同一文書

4

の中でどれが重要文または重要情報であるかは、観点の相違や個人差に左右されるので、必ずしも利用者にとって適切な文や情報が抽出されるとは限らない。また、例えば結論に相当する章または項の文章が抽出されても、その中に必要な文や情報が含まれているとは限らない。さらに、重要文のみを抄録として表示するシステムでは、接続詞などによる文間の関係が壊され、誤解を生むような文等が表示される可能性がある。

【0007】

【発明が解決しようとする課題】前述したように、従来の文書検索システムでは、膨大な量の検索結果から必要な文書や情報を取出すため、利用者には多大な時間と労力が必要となり、負担が大きい。また、文書中から抄録や要約を作成したり、指定の章または項の部分抽出する技術を利用して、検索された文書から抜粋部分を表示するシステムが開発されている。しかし、表示された抜粋部分が、利用者にとって必ずしも適切であるとは限らず、結果的に利用者の必要な抜粋文や情報が得られない場合も多い。

【0008】本発明の第1の目的は、検索された文書の全文から所望の抜粋文や情報のみを抽出して表示できるようにして、利用者の文書検索作業に要する負荷の軽減を図ることにある。

【0009】本発明の第2の目的は、同一文書に対して様々な観点や基準に基づいて抜粋文や情報を抽出して表示できるようにして、多様な検索目的に応じた検索処理を実現することにある。

【0010】本発明の第3の目的は、検索された文書の全文から所定の基準に基づいて抜粋文や情報のみを抽出して表示する場合に、優先度の順序に従って表示する抜粋文や情報の量を調整できる検索処理を実現することにある。

【0011】

【課題を解決するための手段】本発明は、文書を構成する原文データを格納した原文データベースから、例えば意見、提言等のように文章の内容を識別するための複数種類の文の型を設定し、この各文の型に分類した文単位の抜粋文データを作成し、抜粋文データベースとして格納する手段を備えた文書検索システムである。さらに、本システムは、抜粋文データを例えば接続詞を除去したような所定の形式に整形する整形手段と、指定された文の型に対応する抜粋文データを抽出する抽出手段と、抽出された全ての抜粋文データを一覧的に表示する表示手段とを備えている。

【0012】

【作用】本発明では、利用者が予め設定されている複数の文の型から所望の型を選択すると、抽出手段は抜粋文データベースから文の型に対応する抜粋文データを抽出する。例えば意見、提言の文の型を選択すると、検索された文書からその型に分類された全ての抜粋文データが

5

抽出される。表示手段は、抽出された全ての抜粋文データを一覧的に表示する。これにより、利用者は、表示画面上において、例えば意見、提言の文の型に含まれる抜粋文データのみを確認するだけで、必要な文章や情報を得ることができる。また、利用者は文の型を変更するだけで、変更した文の型に含まれる抜粋文データを表示画面上で得ることができる。したがって、文の型を指定するだけで、利用者の必要な文、情報、およびそれらを含む文書を簡単に検索することができることになる。

【 0 0 1 3 】

【 実施例 】 以下図面を参照して本発明の実施例を説明する。

(システムの構成) 図1 は本実施例に係わる文書検索システムの要部を示すブロック図である。

【 0 0 1 4 】 本システムは概略的に、本実施例の文書検索処理に関係する各種処理及びシステム全体の制御を行なう中央処理部 (CPU) 1、表示コントローラ2、表示部3、記憶部4、入力コントローラ10および入力部11からなる。

【 0 0 1 5 】 表示部3 は例えば液晶ディスプレイからなり、表示コントローラ2 の制御により 検索結果や文書等を表示する表示画面を有する。入力部11 は入力コントローラ10 の制御により、検索命令や各種の選択、指定等の入力を行なうキーボード11a とマウス11b を有する。

【 0 0 1 6 】 記憶部4 は、CPU1 の各種処理を行なうためのメインメモリ (RAM)、及びハードディスク装置等のファイル装置からなり、後述する本実施例に係る抜粋文データベース5、原文データベース6、単語インデックス7、型判定用辞書8及び接続詞辞書9を格納している。

(第1の実施例) 前記のシステムにおいて、第1の実施例として図2 に示すような機能ブロック図を示す。図2 において、原文加工部20、検索部21、抜粋部22、インタフェース部23および抜粋インタフェース23a は、CPU1 が実行するプログラムに対応している。

【 0 0 1 7 】 インタフェース部23 は、利用者の検索要求をキーワード、あるいは自然言語からなる文または文章の形式で受け付け、これを検索部21 に渡す。検索部21 は、原文データベース6 及び単語インデックス7 を参照し、該当する文書のリストやその内容を検索結果としてインタフェース部23 に渡す。インタフェース部23 はその検索結果を表示部3 の表示画面上で、利用者に提示する。

【 0 0 1 8 】 インタフェース部23 に組込まれている抜粋インタフェース23a は、検索結果として提示される個々の文書の全文の中から、利用者の指定した文の型に対応する文のみを抜粋して表示画面上に表示する機能を有する。

6

【 0 0 1 9 】 利用者が文の型を選択すると、抜粋インタフェース23a は選択された型を抜粋部22 に渡す。抜粋部22 は、原文データベース6 及び抜粋文データベース5 を参照し、全文の中から該当する型の文のみを抜粋し、文中の接続詞を除去して抜粋インタフェース23a に渡す。抜粋インタフェース23a はこれらの抜粋文を原文と対比させながら表示画面上に表示する。

【 0 0 2 0 】 原文加工部20 は、図3 に示すように、大別して型判定部20a および整形部20b からなり、原文データベース6 から抜粋文データベース5 を作成する。

【 0 0 2 1 】 型判定部20a は、型判定用辞書8 を参照し、与えられた文書の各文がどの型に当てはまるかを判定する。換言すれば、予め設定された複数の型に、文書の各文を分類するための判定処理を実行する。一方、整形部20b は、接続詞辞書9 を参照し、与えられた文書の各文を所定の形式、即ち本実施例では接続詞を除去した形式に整形処理する。

【 0 0 2 2 】 原文加工部20 は、利用者が型の選択をする度にリアルタイムに実行して、抜粋文データベース5 を作成する方式でもよいし、またはシステムの使用前 (検索処理前) にバッチ処理的に抜粋文データベース5 を作成する方式でもよい。本実施例では、高速処理の可能な後者の方式を想定して説明する。

【 0 0 2 3 】 なお、図2 において、破線で描かれた矢印はシステム使用前のデータの流れを意味し、実線で書かれた矢印はシステム使用時のデータの流れを意味する。(原文加工処理) 前記のような原文加工部20 の具体的な動作を図3 乃至図5 を参照して説明する。

【 0 0 2 4 】 まず、原文加工部20 は、原文データベース6 から文書単位に文書を構成する各文 (原データ) を読み込む (ステップS1)。原文データは、図3 に示すように、例えば「 日本的な雇用は、…とされている。」や「 その背景には…と言う 事実がある。」等の各文である。

【 0 0 2 5 】 次に、型判定部20a は型判定用辞書8 を参照し、各文に該当する型を判定し、その型を示すタグを付与する (ステップS2)。ここで、型判定用辞書8 は、例えば図5 (A) に示すように、予め設定された文の型とその型に対応する複数のパターンとを組とする各項目からなる。パターンは、例えばその文の型を的確に表現する文節からなる。

【 0 0 2 6 】 具体例として、図5 (A) に示すように、例えば文の型として「 意見・提言」、「問題提起」、「予想・推定」等が設定されている。「意見・提言」の型には、例えば「 …だと考えられる」や「 …べきではないだろうか」等のパターンがある。この具体例では、パターン「 …であると思う」は「意見・提言」の型と「予想・推定」の型の両方に割り当てられている。即ち、このパターンを含む文は、「意見・提言」型の文を

表現するし、また「予想・推定」型の文も表現する。このように、複数の型に合致する文もあるため、各型は必ずしも相互に排他的ではない。

【0027】また、型判定用辞書8に用意する型の種類は、例えば取り扱う文書の種類に応じて設定する方式でもよい。例えば新聞記事の検索処理では「事実」、「発言」等の型を用意し、技術論文の検索処理では「方法」や「結論」などの型を用意する。

【0028】さらに、整形部20bは接続詞辞書9を参照し、各文を所定の形式、本実施例では接続詞を除去した形式に整形処理する(ステップS3)。具体例として、図3に示すように、例えば「ところが、企業は…変化した。」と言う文から接続詞「ところが、」を除去する。

【0029】原文加工部20は、前記のように加工処理した抜粋文データを抜粋文データベース5に蓄積する(ステップS4)。即ち、原文加工部20は、図3に示すように、原文データ6から判定された型に分類し、接続詞を除去した抜粋文データ5aを作成する。例えば「ところが、企業は…変化した。」という原文は、「…20 した。」というパターンにマッチする「過去」という型のタグが付与され、さらに接続詞「ところが、」を除去された抜粋文に変換される。

【0030】本実施例では、利用者は、後述するように、検索結果の個々の文書からどのような型の文を抜粋して表示するかを指定することができる。この場合、全文の中から「意見・提言」の型の文のみを抜粋表示したり、「結論」の型の文のみを抜粋表示したりしたりして、型を切替えることが可能である。また、図5(B)に示すように、例えば「意見・提言」と「予想・推定」30 という二つの型を同時に選択してもよい。前述したように、文書の全文において、複数の型に関係する文が存在するため、図中のペン図のように、斜線部分の範囲の文が抜粋されるとになる。なお、全文の中で、設定した型に含まれない文も当然ながら存在する。

(抜粋処理)次に、抜粋インタフェース23aによる抜粋文の表示処理について、図6と図7を参照して説明する。

【0031】まず、利用者は、図6(A)に示すように、表示部3の表示画面上に表示された抜粋文の型を選択するためのメニュー画面から、所望の型を入力部11のマウス11bを操作して選択する。ここで、抜粋インタフェース23aは、抜粋部22を介して抜粋文データベース5から文書に対応する抜粋文データを取り出す(ステップS10)。

【0032】抜粋インタフェース23aは、抜粋文データから両者が選択した型のタグが付与された抜粋文の全てを抽出する(ステップS11)。この抽出した抜粋文の文番号を記憶する(ステップS12)。そして、表示画面上に抜粋画面に例えば箇条書き形式で、抽出した全50

ての抜粋文を表示する(ステップS13)。

【0033】具体例として、図6(A)に示すように、型として「意見・提言」が選択された場合に、その型に合致する抜粋文(b, e, h, i)を表示する。

【0034】ここで、表示画面上において、抜粋インタフェース23aは、抜粋画面と原文画面を並べて表示する。抜粋インタフェース23aは、抜粋部22を介して原文データベース6から該当する文書の原文データを取り出す(ステップS14)。そして、原文画面上に取り出した原文データを表示する。このとき、抜粋インタフェース23aは、記憶した抜粋文の文番号に対応する原文(b, e, h, i)を例えばカラー表示やアンダーライン表示等による強調表示する(ステップS15)。

【0035】即ち、原文画面上において、「意見・提言」型に合致して、接続詞「しかし、」を除去した抜粋文bに対応する文番号2の原文が強調表示される。また、抜粋文画面上に抜粋文iが表示されているが、これは原文画面中に収まりきらなかった文のうち「意見・提言」の型に該当した文である。

【0036】ここで、同一文書において、図6(B)に示すように、利用者が型を「予想・推定」に切替えたとき、原文画面中の文番号4, 5の原文は「予想・推定」の型に該当するため抜粋されている。この場合、文番号5に対応する抜粋文データeは、「意見・提言」の型と「予想・推定」の型の両方に該当している。

(パターン表示処理)以上のように、利用者が抜粋文の型を選択するには、選択肢として表示されている各型の意味を利用者にわかりやすく表示することが望ましい。例えば、「意見・提言」という型を選択すると、実際には「…だと考えられる」や「…べきではないだろうか」のようなパターンの文が抜粋されるということを利用者が知っていれば、所望の情報を得やすいと考えられる。

【0037】そこで、図8(A)に示すように、利用者が型を選択すると、抜粋部22が型判定用辞書8から選択した型に対応する全てのパターンを取出して抜粋インタフェース23aに出力する(ステップS20, S21)。抜粋インタフェース23aは、同図(B)に示すように、例えば選択した「意見・提言」の型に対応するパターンを列挙した一覧表示を行なう(ステップS22)。

【0038】この機能により、利用者は、どの型を選択すると、どのような表現を含む文が抜粋されるかを予め知ることができる。

(選択表示処理)これまでの説明では、抜粋文データにおいて選択中の型のタグがつけられている文はすべて抜粋文として表示するものとしていた。この変形例として、選択中の型のタグがついている文をとりだした後、この中からある尺度に基づいて一つ以上の重要文を選択し、それを表示する方式を説明する。

【0039】まず、図9に示すように、選択した型のタ

グが付与された抜粋文の全てを抽出するまでの処理は、これまでの処理と同様である(ステップS30, S31)。

【0040】次に、検索要求に含まれる検索キーワードを最も多く含む文が重要文であるとし、この尺度に基づいて該当する抜粋文を抜粋画面上に箇条書き形式で表示する(ステップS32, S33)。

【0041】具体例として、図10に示すように、選択した例えば型「目的・ねらい」に対して、全文6文の中で、該当する文である文番号2, 4の各文が抜粋文として抽出される。この2文の中で、文番号2の文は「分散型」と「アーキテクチャ」という2つの検索キーワードを含むが、文番号4の文はひとつも含んでいない。したがって、文番号2の文の方が文番号4よりも重要な文であるとし、この文番号2の文のみを抜粋文として表示する。

【0042】これにより、利用者は、選択した型と共に、検索要求の検索キーワードを含む最も適切な抜粋文を参照できる可能性が高くなる。

(特定文表示処理)これまでの説明では、利用者は型の選択権のみをもち、ある型が選択された場合、その型に対応する全てのパターンにマッチする文が抜粋されることになっていた。この変形例として、特定のパターンにマッチする文のみを抜粋するように利用者に指定させる方法を説明する。

【0043】まず、図11に示すように、選択した型に対応するパターンを、型判定辞書8から取り出すまでの処理は、これまでのパターン表示処理と同様である(ステップS40, S41)。

【0044】ここでは、図12(A)に示すように、例えば「意見・提言」の型に該当するパターンを取り出している。このとき、抜粋文画面上に各パターンと共に、各パターンにマッチする文の数を表示する(ステップS42)。この例では、「…だと考えられる」のパターンの文は5文、「…であると思う」のパターンの文は3文であることが表示されているので、この二つのパターンを選択すると抜粋文画面には合計8文が表示されることが事前にわかる。同様に、各パターンにマッチした文の長さや行数などの情報を表示すれば、利用者は画面にちょうど収まるくらいの分量の文を抜粋表示させることもできる。

【0045】この表示されたパターンの中から、利用者が例えば「…だと考えられる」と「…であると思う」の各パターンを選択すると、抜粋部22は抜粋文データベース5から該当する抜粋文データを取り出して抜粋インタフェース23aに返す(ステップS44)。

【0046】抜粋インタフェース23aは、図12(A)に示すように、利用者が選択した型に該当し、かつ選択した特定パターンにマッチする抜粋文のみを抜粋文画面に表示する(ステップS45)。この場合、前述

のパターンにマッチした抜粋文の中から重要文のみを表示する処理を適用してもよい。

【0047】以上のように、型の選択に加えてパターンの選択も利用者に行わせると、利用者は抜粋文の分量をある程度制御することができる。

(抜粋文の分量制御)次に、前記の変形例として、型判定辞書8に用意したパターンに予め優先順位を設定し、パターンの優先順位に基づいて表示する文の分量を制御する方式について説明する。

【0048】この変形例では、図14(A)に示すように、型判定辞書8の構造として、型とパターンとの項目以外に、パターン間の優先順位を示す情報が含まれている。具体例として、「意見・提言」という型に対応するパターンの中では、「…だと考えられる」というパターンの優先順位が一番高いことが示されている。

【0049】ここで、現在選択中の文書について、「意見・提言」の型に対応する各パターンにマッチする文の数をカウントすると、図14(B)に示すような文数情報が得られると想定する。即ち、「…だと考えられる」と「…べきではないだろうか」のパターンの文のみ表示すれば合計3文、さらに次に優先度の高い「…を提案する」のパターンの文まで表示すれば合計6文が表示されることになる。

【0050】この変形例では利用者が、図13に示すように、表示画面上において抜粋表示する文の数(表示文数)nを指定する(ステップS50)。ここでは、例えば図14(C)に示すように、利用者が表示文数nを5文以内に制限したと想定する。

【0051】抜粋部22は優先順位が記述されている型判定辞書8を参照し、利用者の指定した分量nと優先順位までのパターンにマッチする文の分量mとを比較する(ステップS51, S52)。そして、優先順位の高いパターンにマッチする抜粋文を、指定分量nの範囲内で段階的に選択していく(ステップS54~S55)。最終的に、指定分量nの範囲内でかつ優先順位の高いパターンにマッチする抜粋文を抜粋インタフェース23aに出力する。

【0052】抜粋インタフェース23aは、図14(C)に示すように、利用者が指定した表示文数(5文)の範囲内の抜粋文(ここでは3文)を抜粋画面に表示する(ステップS56)。ここでは、優先順位が1位と2の「…だと考えられる」と「…べきではないだろうか」の各パターンにマッチする抜粋文のみを表示する。

【0053】なお、ここでは文の分量の例として文数で説明したが、文の長さ、行数などをもとに表示の制御をすることも可能である。また、利用者に文の分量を指定させるかわりに、優先順位何番目までのパターンにマッチする文を表示させるかを指定させてもよい。

(型判定辞書の更新処理)これまでの説明では、型判定辞書8は事前に用意された固定的なものを想定してい

た。この変形例では、型判定辞書8に新たなパターンを追加登録したり、または逆に不適切と思われるパターンを削除する方式を説明する。

【0054】まず、利用者の入力に応じて、型判定辞書8の更新用画面を表示する(ステップS60)。次に、更新対象の型とパターンを入力する(ステップS61)。具体例としては、例えば図17(A)に示すように、利用者が原文画面上で、パターンとして登録したい範囲(ここでは、すべきだ)をマウス11bにより指定する。そして、同図(B)に示すように、登録対象の型を指定する。ここでは、「意見・提言」の型に「…すべきだ」という表現のパターンを追加する場合である。

【0055】これらの入力に応じて、抜粋部22は型判定辞書8の該当する型に関するエントリを検索する(ステップS62)。そして、該当する型のエントリとして、入力した追加パターンを型判定辞書8に登録して更新する(ステップS63)。

【0056】具体例として、図16(B)に示すように、「意見・提言」の型のパターンのひとつとして、「…すべきだ」という表現を追加する。

【0057】一方、追加だけでなく、逆に指定のパターンを削除する更新処理もある。具体例として、図16(A)に示すように、例えば、「意見・提言」の型に対応するエントリとして「…べきではないだろうか」という記述があった場合、利用者がこれを好ましくないと思う場合もあると考えられる。このような場合には、更新処理として削除を指定し、「…べきではないだろうか」という表現のパターンを削除する。

【0058】このような機能により、利用者にとって、予め設定されている型判定辞書8に不備があり、用意されている型に対応するパターンの中に必要と思われるパターンが登録されていない場合には、そのパターンを登録することができる。また、逆に型の表現としては不適切と思われるパターンが存在する場合には、そのパターンを削除することができる。型判定辞書の各エントリは、型とパターンの対という単純な構造をしているので、利用者がこれらの追加あるいは削除を行うことは容易である。したがって、利用者にとって有効な型判定辞書8を再構築することが可能となる。

(優先度修正処理) 前述の分量制御処理において、型判定辞書8のパターン間に優先順位を付加する方式について説明したが(図13と図14を参照)、この優先順位を利用者に適応するように修正する方式について説明する。

【0059】まず、図19(B)に示すように、選択した型に対応する抜粋文データを抜粋画面上に表示する(ステップS70)。ここで、抜粋文データを選択するための型判定辞書8には、図19(A)に示すように、各パターンには優先順位を示す情報が付加されている。

【0060】次に、利用者がマウス11bにより、抜粋

画面上において重要だと思われる抜粋文を指定する(ステップS71)。ここでは、図19(B)に示すように、具体例として、「…だと考えられる」というパターンの文と「…すべきではないか」というパターンの文が表示されているが、利用者は、内容的にみて「…すべきではないか」の文が重要であると判断し、マウス11bにより指定したと想定する。

【0061】抜粋部22は指定された抜粋文データから、その型に含まれるパターンを抽出し、型判定用辞書8においてそのパターンの優先順位を高い方に修正する処理を行なう(ステップS72, S73)。

【0062】具体例として、図19(A)に示すように、「意見・提言」の型のパターンとして、当初の型判定用辞書8では、「…だと考えられる」のパターンの優先順位が第1位で、「…べきではないだろうか」のパターンの優先順位が第2位である。ここで利用者の「…べきではないだろうか」のパターンの優先度の修正を指定したことにより、図19(A)に示すように、両者の優先順位が逆転するように、型判定用辞書8の内容が更新される。

【0063】これにより、以後の分量制御処理を伴う抜粋処理において、更新した型判定用辞書8を使用することにより、利用者に適応した抜粋文を抽出して表示することになる。

(複数文書の表示処理) これまでの説明では、検索結果の中の一文書を選択した時の処理が中心であった。この変形例では、検索結果の中の複数の文書に対して同一の型やパターンにより、抜粋処理を行なう方式について説明する。

【0064】ここで、一般的に、検索結果は、検索要求に近いものから順に提示されることが多いが、選択した型やパターンに該当する抜粋文が多い順に文書を提示する方法も考えられる。このようにすれば、あとの方に提示される文書ほど抜粋される情報が少ないことがわかるので、複数の文書から効率的に情報を得ることが可能となる。

【0065】まず、図20に示すように、利用者の入力に応じて型とパターンが選択されると、検索されている複数の文書(ここではA, B, Cの3文書)をソートする処理を行なう(ステップS80~S82)。そして、ソートした順序で各文書を表示する(ステップS83)。

【0066】ここで、図21(A)に示すように、検索された各文書A, B, Cの抜粋文データには、各型やパターンに該当する文がいくつあったかという情報が記述されている。例えば、文書Aの抜粋文データには、「意見・提言」の型に該当する文が10文、「問題提起」の型に該当する文が0文、「予想・推定」の型に該当する文が4文含まれることなどが記されている。

【0067】利用者の選択している型が「意見・提言」

である場合、この型に該当する文の数はA、B、Cの順で多いので、図21(B)に示すように、この順序で文書のリストを表示する。利用者の選択している型が「予想・推定」である場合、この型に該当する文の数はB、C、Aの順で多いので、図21(C)に示すように、この順序でリストを表示する。

(第2の実施例) 次に、本発明の第2の実施例を、図22に示す機能ブロック図を参照して説明する。本実施例は、前述した第1の実施例との相違点は、抜粋インタフェース23bが原文表示専用の画面をもたないということである。以下、第1の実施例と異なる点について説明する。したがって、異なる点以外の機能は第1の実施例と同様である。

(全文選択処理) 本実施例は、図24に示すように、予め設定される複数の抜粋文の型の中に、「全文」という項目を含む。この項目を選択すると、抜粋画面上に、加工処理されていない原文の全文を直接表示する。

【0068】即ち、図23のフローチャートに示すように、両者が型を選択する場合に、「全文」の項目を選択すると、抜粋部22は原文データベース5から原文データからなる全文を取り出し、抜粋インタフェース23bに出力する(ステップS90、S91、S95)。抜粋インタフェース23bは、その原文データを抜粋文データとして抜粋画面上に表示する(ステップS96)。

【0069】一方、「全文」以外の型を選択した場合には、前述の第1の実施例と同様に、抜粋文データベース5から選択した型のタグが付与された抜粋文を全て取り出して、抜粋画面上に箇条書き形式で表示する(ステップS91のNO、S92～S94)。具体例として、図24に示すように、「意見・提言」の型が選択されると、その型に対応する抜粋画面上に表示する。

【0070】本実施例によれば、原文画面を用意することなく、「全文」という型を用意するだけで、原文データを抜粋画面上で直接参照することが可能となる。したがって、利用者は、必要に応じて抜粋文データまたは原文データを簡単に選択して参照することが可能である。

(第3の実施例) 次に、本発明の第3の実施例を、図25に示す機能ブロック図を参照して説明する。本実施例の特徴は、文の型を判定するためのパターンとして、前述の第1と第2の実施例とは異なり、言語解析処理により得られる言語情報を利用する点にある。

【0071】即ち、前述の第1と第2の実施例では、文の型を判定するためのパターンは、例えば「…だと考えられる」のような表層文字列である。即ち、「…だと考えられる」のような表層文字列を含む抜粋文は、例えば「意見・提言」の型に相当すると判定する。

【0072】これに対して、本実施例は、原文の言語解析により得られる形態素情報(品詞、活用形等)や構文情報(主語、述語、係り受け等)の高度な言語情報を利用したパターンマッチング処理を実行して、最適な抜粋

文を表示する。

【0073】以下、前述の第1と第2の実施例との相違点のみについて説明する。したがって、異なる点以外の機能は第1と第2の実施例と同様である。

(原文加工処理) 図26のフローチャートに示すように、まず、原文加工部20は、原文データベース6から文書を構成する各文(原文データ)を読み込む(ステップS100)。原文データは、図27に示すように、例えば「日本的な雇用は、…である。」や「おそらく、その背景には…であろう。」等の各文である。

【0074】ここで、図27に示すように、原文加工部20の解析部20cは、予め用意された解析辞書24を参照して言語解析処理を実行し、形態素情報や構文情報等の言語情報を抽出する(ステップS101)。

【0075】型判定部20aは型判定用辞書8を参照し、各文に該当する型を判定し、その型を示すタグを付与する(ステップS102)。そして、接続詞辞書9を参照した整形処理の後に、抜粋文データを抜粋文データベース5に格納する(ステップS103、S104)。このとき、抜粋文データには、型のタグと共に形態素情報、構文情報などの言語情報が付与されている。

【0076】具体例として、図27に示すように、例えば「日本的な雇用は、…である。」という抜粋文データには、「事実」の型を示すタグと形態素情報として「断定の助動詞」の言語情報が付与されている。

【0077】このような言語情報を利用して、利用者が選択した型に対応するパターンマッチング処理を行なう場合に、型判定用辞書8には各型に対応するパターンとして前記の言語情報が用意されている。

【0078】具体的には、図28に示すように、パターンとして表層文字列を列挙する代わりに、例えば言語情報として品詞の種類に関する情報が設定されている。さらに、この言語情報と表層文字列を組合わせたパターンを用意するのが望ましい。

(第4の実施例) さらに、本発明の第4の実施例を、図29に示す機能ブロック図を参照して説明する。本実施例の特徴は、前述の第1乃至第3の実施例とは異なり、型の判定処理と抜粋処理を文単位に限定せずに、文の様々な構成要素を単位とする点にある。即ち、句や節、特定の意味内容を表す表現など、文の様々な構成要素を抽出して表示する。

【0079】本実施例では、抜粋処理により得られるデータを抜粋情報データと称し、原文加工部20の加工処理により得られた抜粋情報データは抜粋情報データベース25に格納される。

(原文加工処理) 以下、図31と図32を参照して本実施例の原文加工処理について説明する。

【0080】まず、原文加工部20は、原文データベース6から文書を構成する各文を読み込む(ステップS110)。原文データは、例えば図32(B)に示すよう

10

20

30

40

50

に、「…首相は、昨夜の首相官邸における…」等の文である。

【0081】ここで、原文加工部20の解析部20cは、予め用意された解析辞書24を参照して既存の言語解析処理を実行し、重文や複文から単文を抽出したり、主語や固有名詞、時間などを表現形態等の構成要素（構成単位）の表現を抽出する（ステップS111）。

【0082】型判定部20aは型判定用辞書8を参照し、各構成単位に対して型を判定し、その型を示すタグを付与する（ステップS112）。整形部20bは、まず接続詞辞書9に記述されている接続詞や接続助詞などを文や節から除去し、次に時間を表す型「when」にマッチした文の構成単位の正規化などを行う（ステップS113）。そして、原文加工部20は型や言語情報を付与された抜粋文と共に、時間や主語等の構成単位を出力する（ステップS114）。

【0083】具体例として、図30に示すような抜粋情報データが作成されて、抜粋情報データベース25に格納される。この具体例では、第一文の前半である「日本の政治は、…であるが、」という部分は、「事実」という型に該当し、一方、後半の「この背景には、…という事実がある。」という部分は、「背景」という型に該当している。このようにひとつの文に複数の種類の情報が含まれる場合にも対処できる。また、前半の「日本の政治は、…であるが、」という節末尾には接続助詞や読点がついているが、整形部はこれらを除去して「…である」という終止形にしている。

【0084】さらに、図30に示すように、第2文からは、「when」という時間に関する情報と「who」という動作主体を表現する情報が抽出されている。ここで、原文中の「平成6年7月」という表現が、整形部20bにより「1994. 7」という形式に正規化されている。

【0085】本実施例では、図32（D）に示すように、型判定用辞書8には例えば5 W1 Hのような抜粋情報の型とそれに対応するパターンが用意されている。この場合、図28に示すように、品詞情報などの言語情報を用いてパターンを記述してもよい。

【0086】このような型判定用辞書8を利用して、図32（A）に示すように、予め用意された5 W1 Hの抜粋情報の型から、利用者が例えば「who」という動作主体を表現する型が選択されると、原文から例えば「首相」、「前首相」、「科学技術庁長官」等のなんらかの動作の主体となっている人物の名前が抽出されて抜粋情報画面に表示される。同様に、利用者が例えば「where」という場所を表現する型が選択されると、同一原文中から例えば「首相官邸」、「国内の原子力発電所」等の場所を表現する言葉が抽出されて抜粋情報画面に表示される。

【0087】

【発明の効果】以上詳述したように本発明によれば、フルテキストの文書を検索する文書検索システムにおいて、第1に、検索された文書の全文から所望の抜粋文や情報のみを抽出して表示できるようにして、利用者の文書検索作業に要する負荷の軽減化を図ることができる。第2に、同一文書に対して様々な観点や基準に基づいて抜粋文や情報を抽出して表示できるようにして、多様な検索目的に応じた検索処理を実現することができる。第3に、検索された文書の全文から所定の基準に基づいて抜粋文や情報のみを抽出して表示する場合に、優先度の順序に従って表示する抜粋文や情報の量を調整できる検索処理を実現することができる。

【0088】換言すれば、利用者の望む型の文や文の構成要素のみを抜粋して表示することにより、大量の検索結果のブラウジングを効率化することができる。さらに、抜粋情報のみをコンパクトに表示することにより、文書全体を直観的にかつ総括的にとらえることを可能にし、画面のスクロールなどに関する労力及び負荷を軽減することができる。さらに、原文中でどの文が重要かは観点によっても変わるし個人差もあるが、抜粋情報の型の切替により、様々な観点から重要な箇所を抜粋することが可能である。

【図面の簡単な説明】

【図1】本発明の実施例に係わる文書検索システムの要部を示すブロック図。

【図2】第1の実施例に係わる機能ブロック図。

【図3】第1の実施例の動作を説明するための概念図。

【図4】第1の実施例の動作を説明するためのフローチャート。

【図5】第1の実施例の動作を説明するための概念図。

【図6】第1の実施例の動作を説明するための概念図。

【図7】第1の実施例の抜粋処理に関する動作を説明するためのフローチャート。

【図8】第1の実施例のパターン処理に関する動作を説明するための図。

【図9】第1の実施例の選択表示処理に関する動作を説明するためのフローチャート。

【図10】第1の実施例の選択表示処理に関する動作を説明するための図。

【図11】第1の実施例の特定文表示処理に関する動作を説明するためのフローチャート。

【図12】第1の実施例の分量制御に関する動作を説明するための図。

【図13】第1の実施例の分量制御に関する動作を説明するためのフローチャート。

【図14】第1の実施例の分量制御に関する動作を説明するための図。

【図15】第1の実施例の更新処理に関する動作を説明するためのフローチャート。

【図16】第1の実施例の更新処理に関する動作を説明

するための図。

【図17】第1の実施例の更新処理に関する動作を説明するための図。

【図18】第1の実施例の優先度修正処理に関する動作を説明するためのフローチャート。

【図19】第1の実施例の優先度修正処理に関する動作を説明するための図。

【図20】第1の実施例の複数文書の検索処理に関する動作を説明するためのフローチャート。

【図21】第1の実施例の複数文書の検索処理に関する動作を説明するための図。

【図2.2】第2の実施例に係わる機能ブロック図。

【図23】第2の実施例の動作を説明するためのフローチャート。

【図24】第2の実施例の動作を説明するための図。

【図25】第3の実施例に係わる機能ブロック図。

【図26】第3の実施例の動作を説明するためのフロー

チャート。

【図27】第3の実施例の動作を説明するための図。

【図28】第3の実施例の型判定用辞書の構造を説明するための図。

【図29】第4の実施例に係わる機能ブロック図。

【図30】第4の実施例の動作を説明するための図。

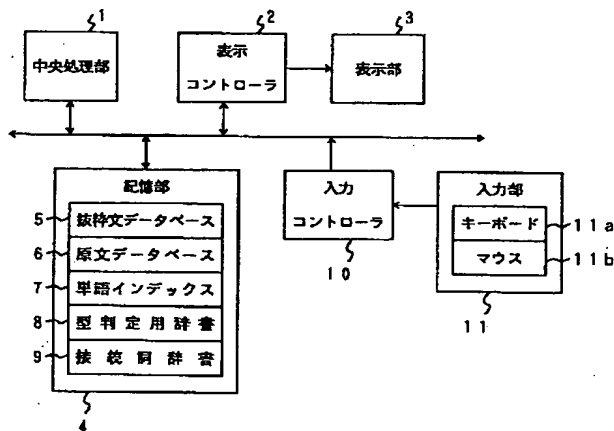
【図31】第4の実施例の動作を説明するためのフローチャート。

【図32】第4の実施例の動作を説明するための図。

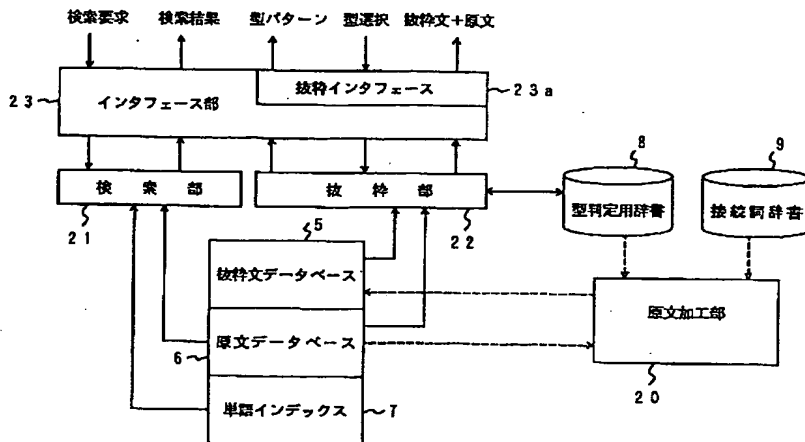
【符号の説明】

1…中央処理部(CPU)、2…表示コントローラ、3…表示部、4…記憶部、5…抜粋文データベース、6…原文データベース、7…単語インデックス、8…型判定用辞書、9…接続詞辞書、10…入力コントローラ、11…入力部、20…原文加工部、21…検索部、22…抜粋部、23…インターフェース部、23a…抜粋インターフェース。

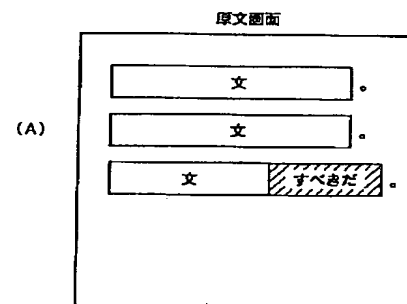
【図1】



【図2】



【図17】



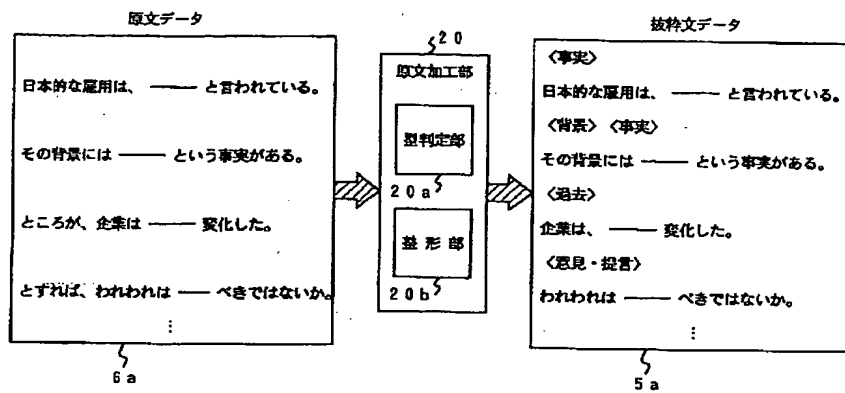
【辞書登録画面】登録語:「…すべきた」
どの型のパターンとして登録しますか?

☒ 意見・提言
☐ 問題提起
☐ 予想・推定
...

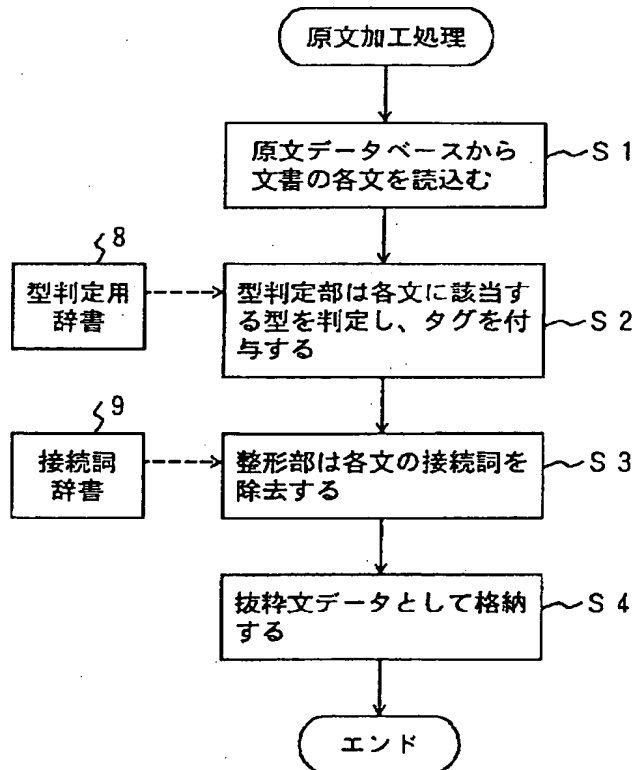
登録 中止

(B)

【 図3 】



【 図4 】

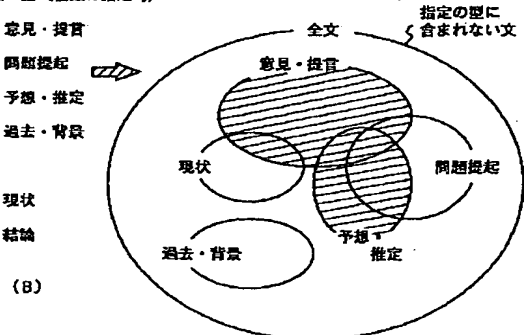


【 図5 】

〈型〉	〈パターン〉
意見・提言	「…だと考えられる」 「…べきではないだろうか」 「…を提案する」 「…であると思う」 ⋮
問題提起	「…問題がある」 「…について考えてみたい」 ⋮
予想・推定	「おそらく、…」 「…なるであろう」 「…かも知れない」 「…であると思う」 ⋮
⋮	⋮

抜特文の型（複数の指定可）

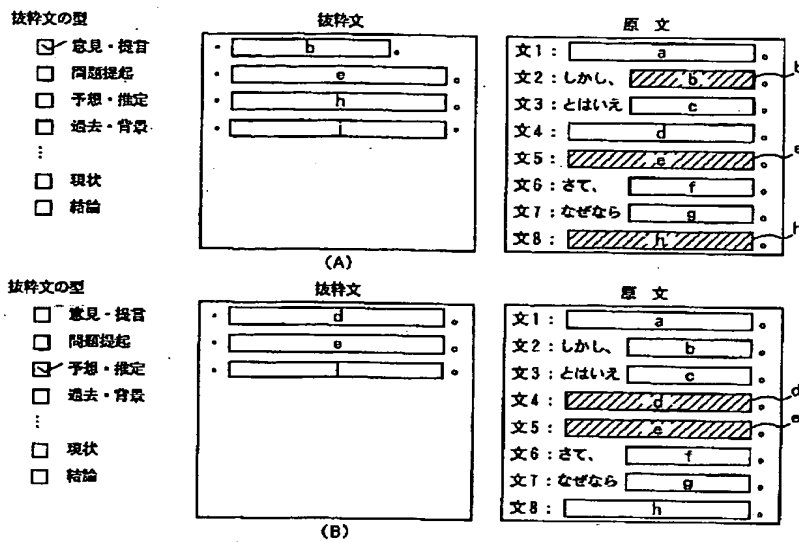
- ☒ 意見・提言
☐ 問題提起
☒ 予想・推定
☐ 過去・背景
:
☐ 現状
☐ 結論



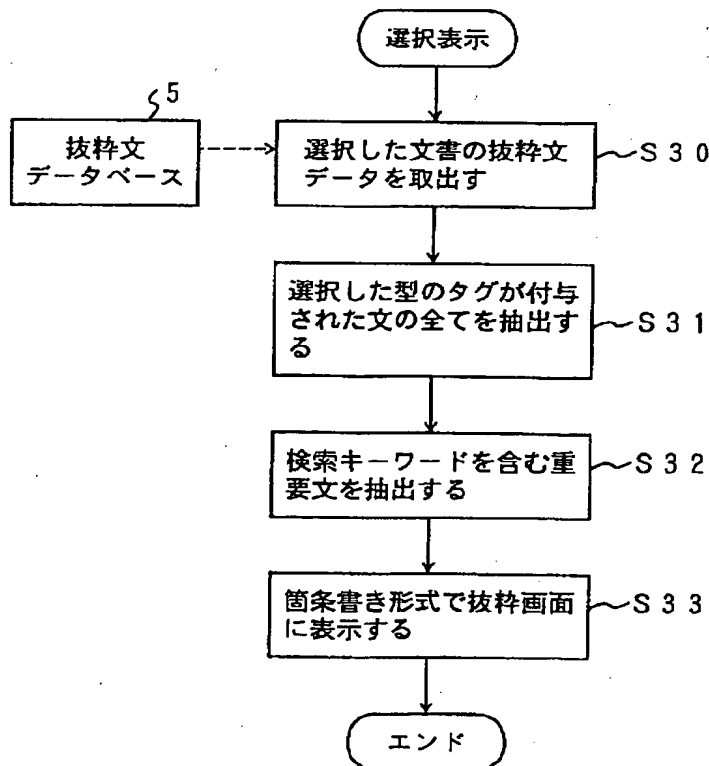
【图28】

〈型〉	〈パターン〉
事実	文末に断定の助動詞がある
予想・推定	文章に推定の副詞がある 文末に推定の助動詞がある

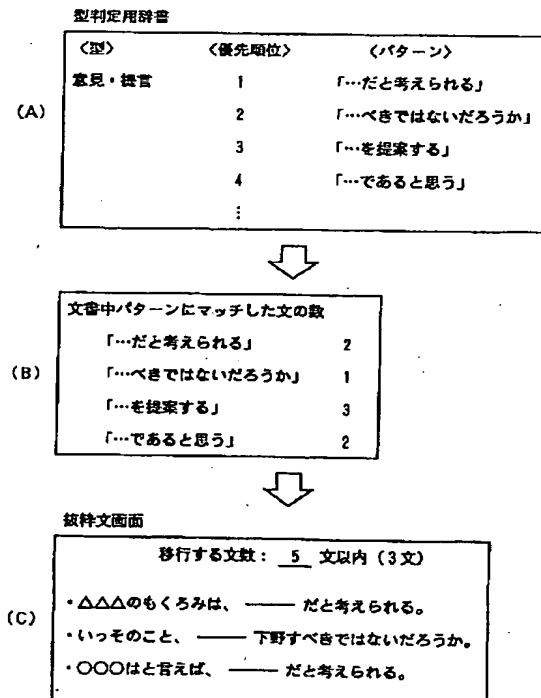
【 図6 】



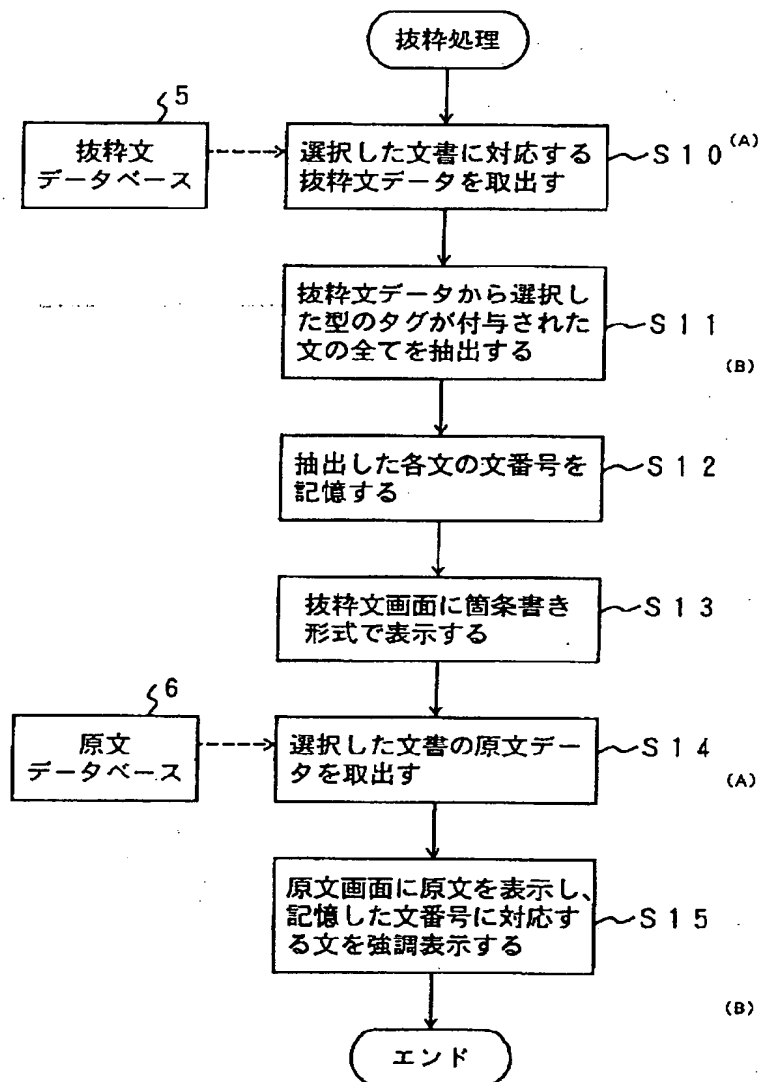
【 図9 】



【 図14 】



【 図7 】



【 図16 】

型: 意見・提言

追加ボタン: すべきだ ▲

パターン: ...だと考えられる
...べきではないだろうか
...を提案する
...であると思う

型判定用辞書の更新

【 図19 】

〈型〉	〈優先順位〉	〈パターン〉
意見・提言	1	「...だと考えられる」
	2	「...べきではないだろうか」
	3	「...を提案する」
	4	「...であると思う」
	⋮	⋮

・△△△のもくろみは、—— だと考えられる。

・いっそのこと、—— 下野すべきではないだろうか。

・〇〇〇はと言えば、—— だと考えられる。

【 図24 】

抜粋文の型

- ☒ 意見・提言
☐ 問題提起
☐ 予想・推定
☐ 過去・背景
 ⋮
☐ 現状
☐ 結論
☐ 全文

抜粋文

・ [文] 。

・ [文] 。

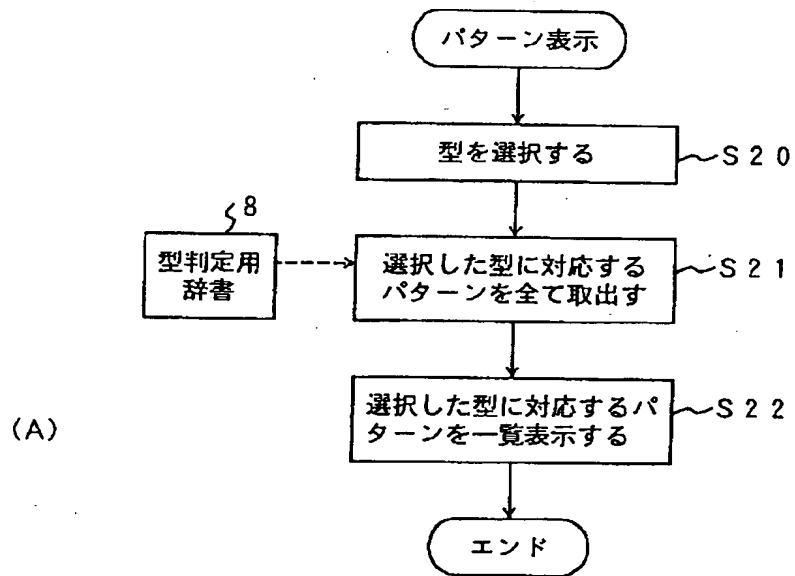
・ [文] 。

・ [文] 。

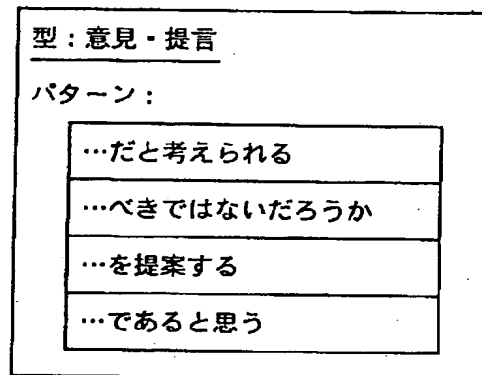
(C)

〈型〉	〈優先順位〉	〈パターン〉
意見・提言	2	「...だと考えられる」
	1	「...べきではないだろうか」
	3	「...を提案する」
	4	「...であると思う」
	⋮	⋮

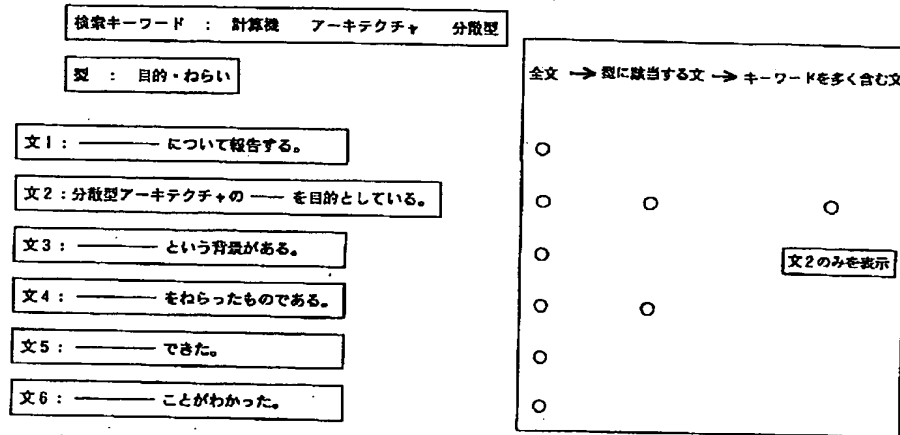
【 図8 】



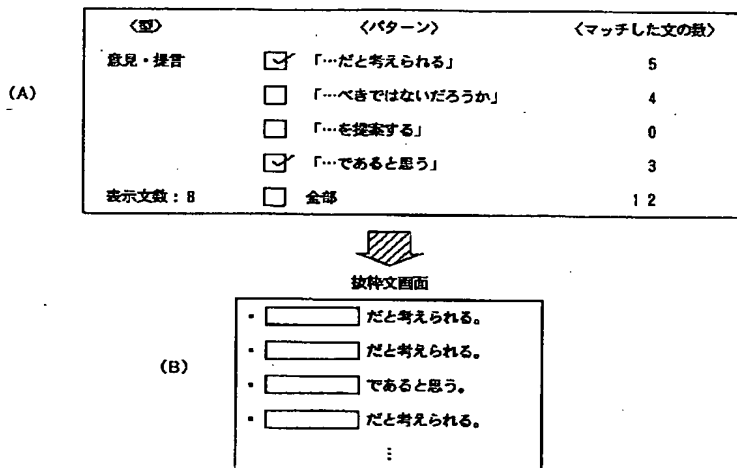
(B)



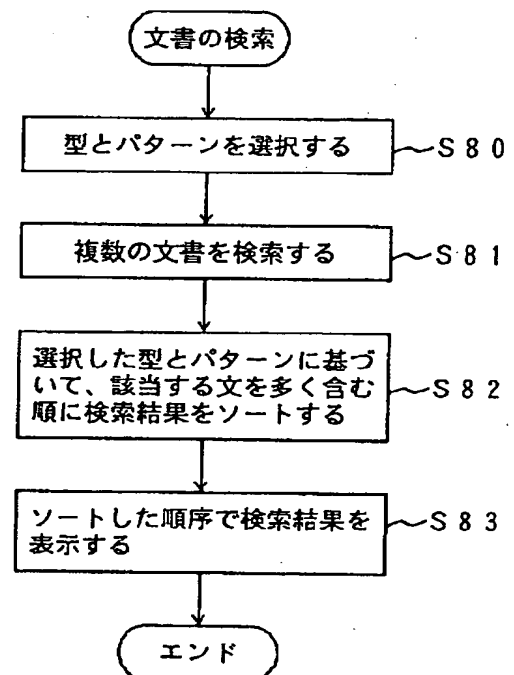
【 図10 】



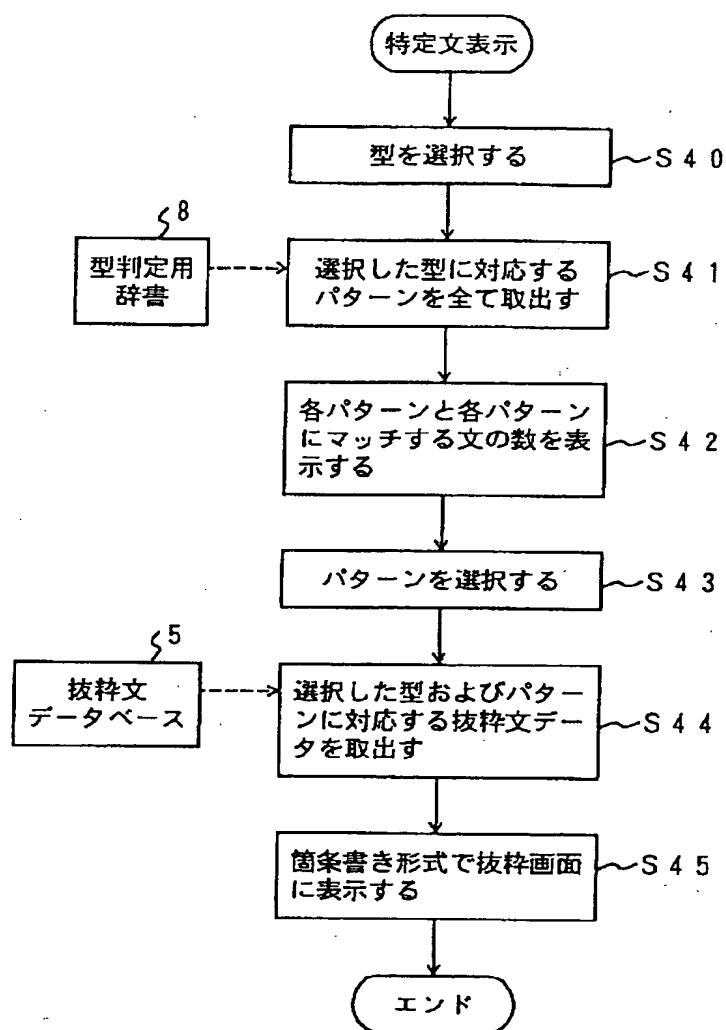
【 図12 】



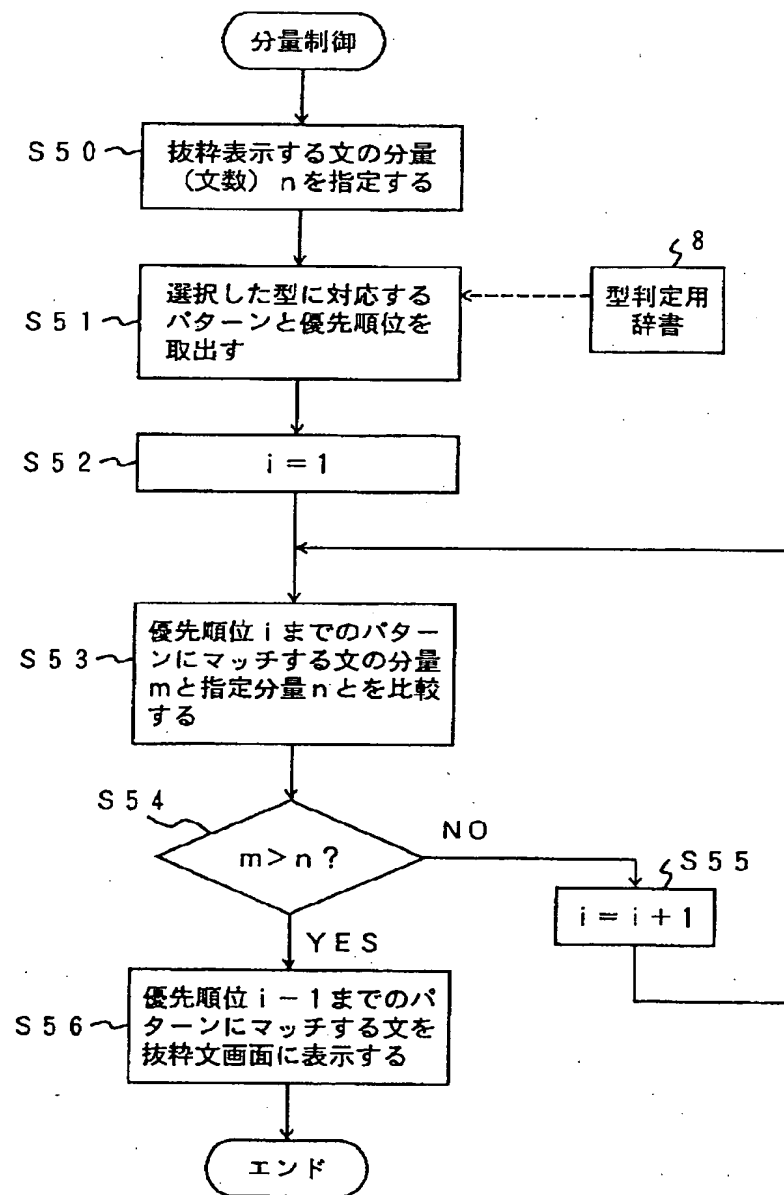
【 図20 】



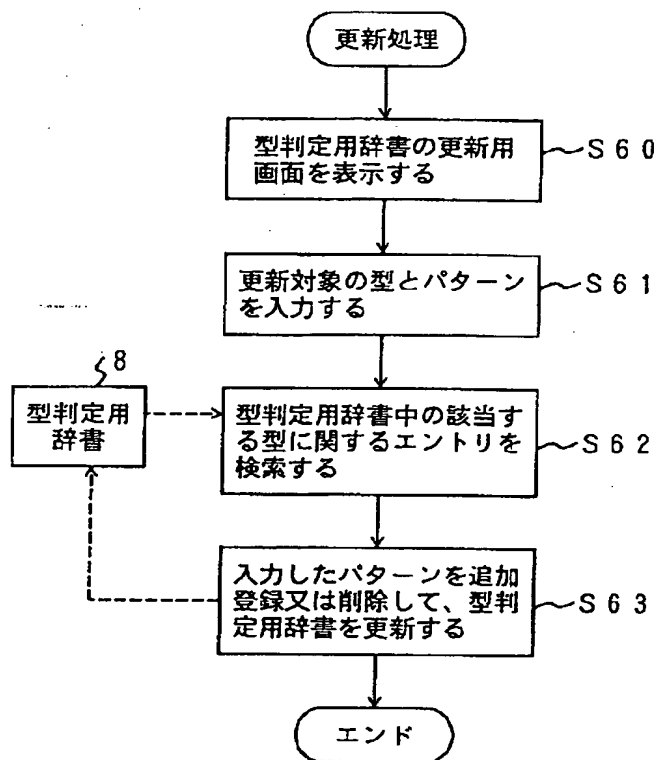
【 図1 1 】



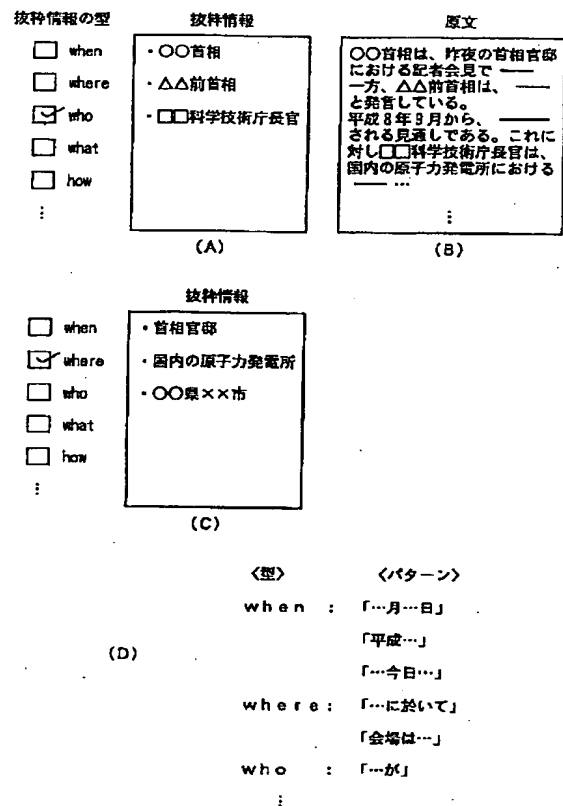
【 図13 】



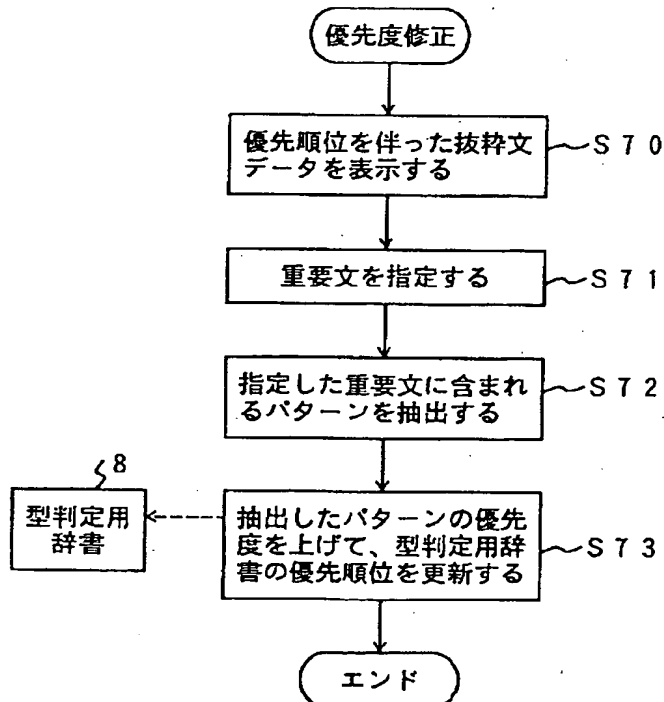
【 図15 】



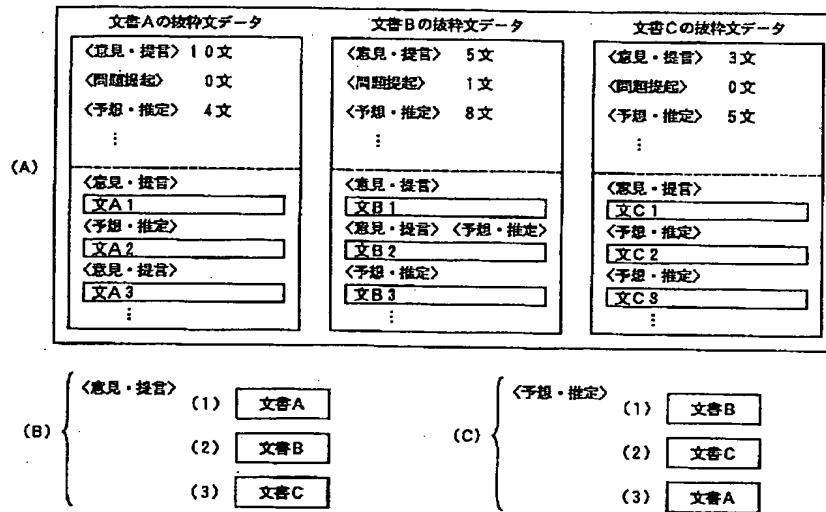
【 図32 】



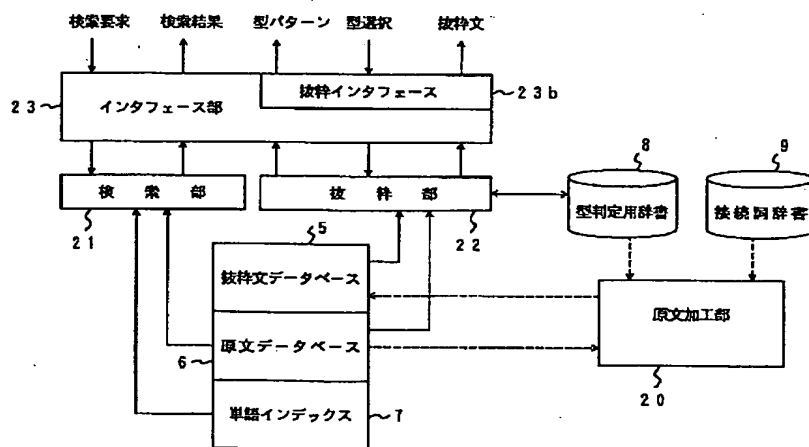
【 図18 】



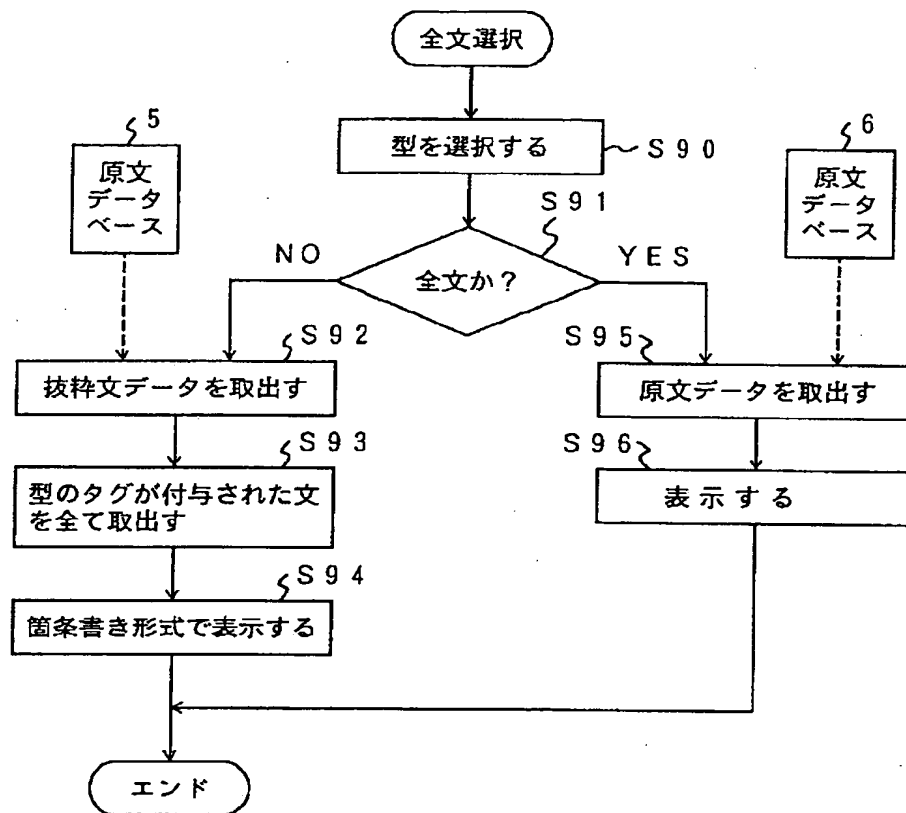
【 図2 1 】



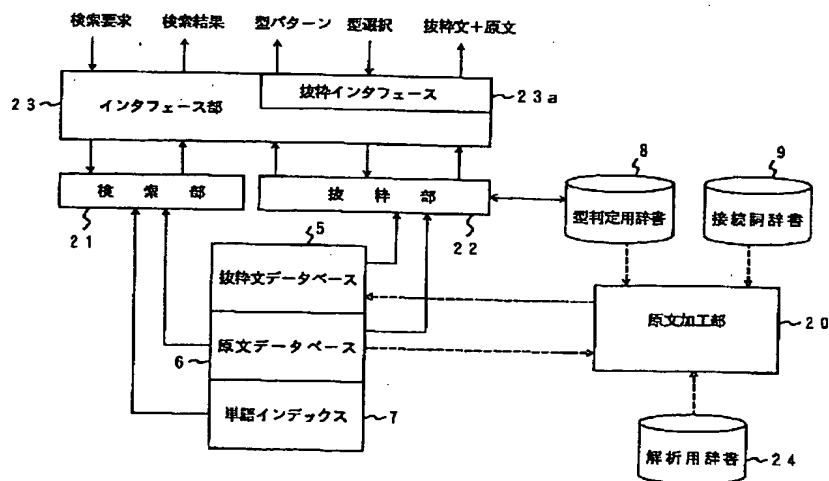
【 図2 2 】



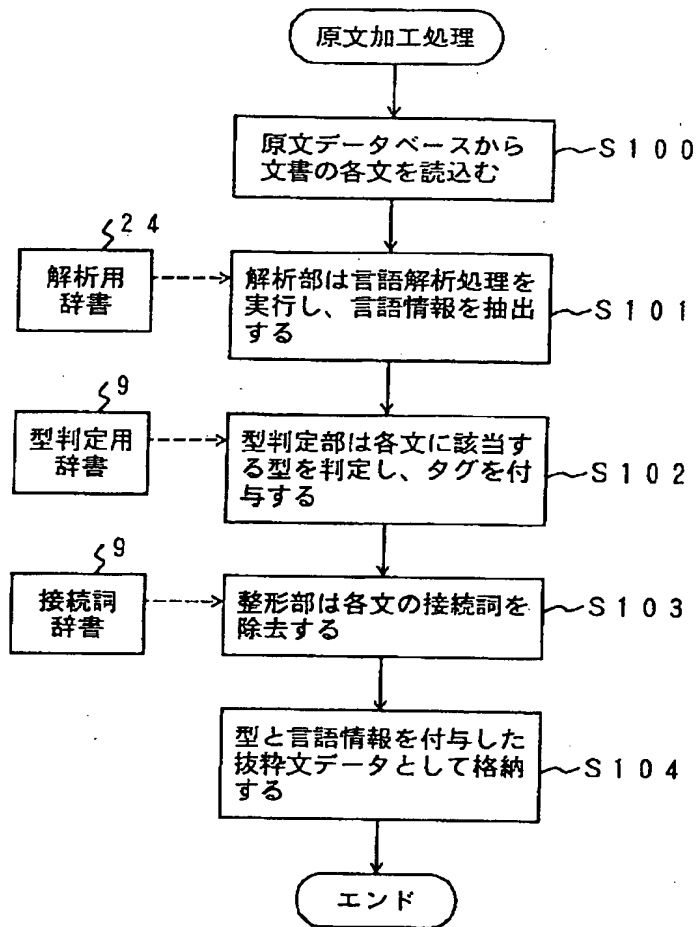
【 図23 】



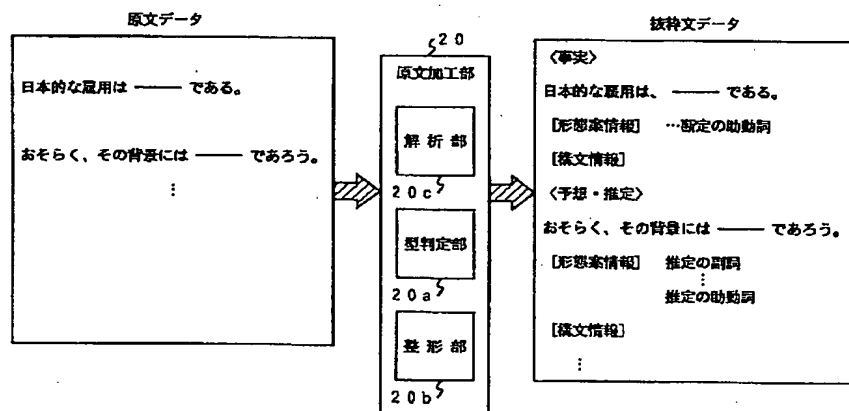
【 図25 】



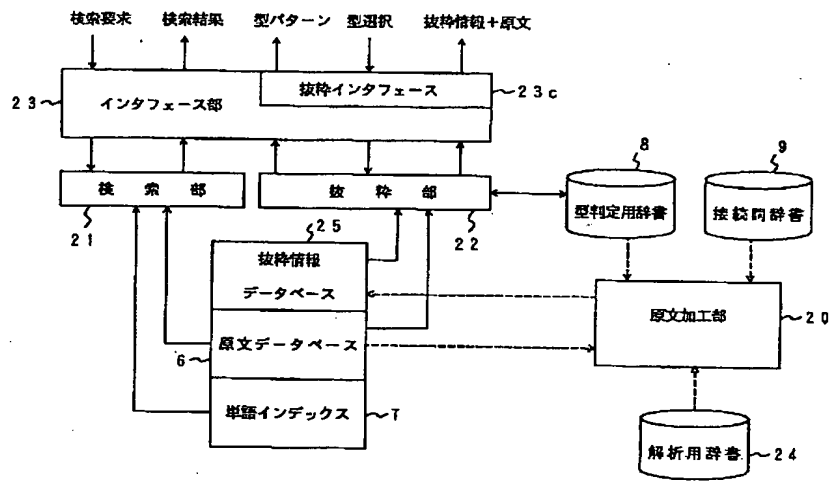
【 図26 】



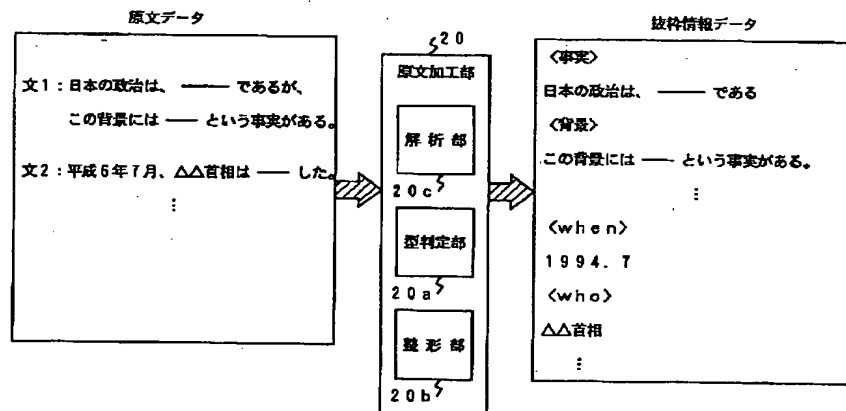
【 図27 】



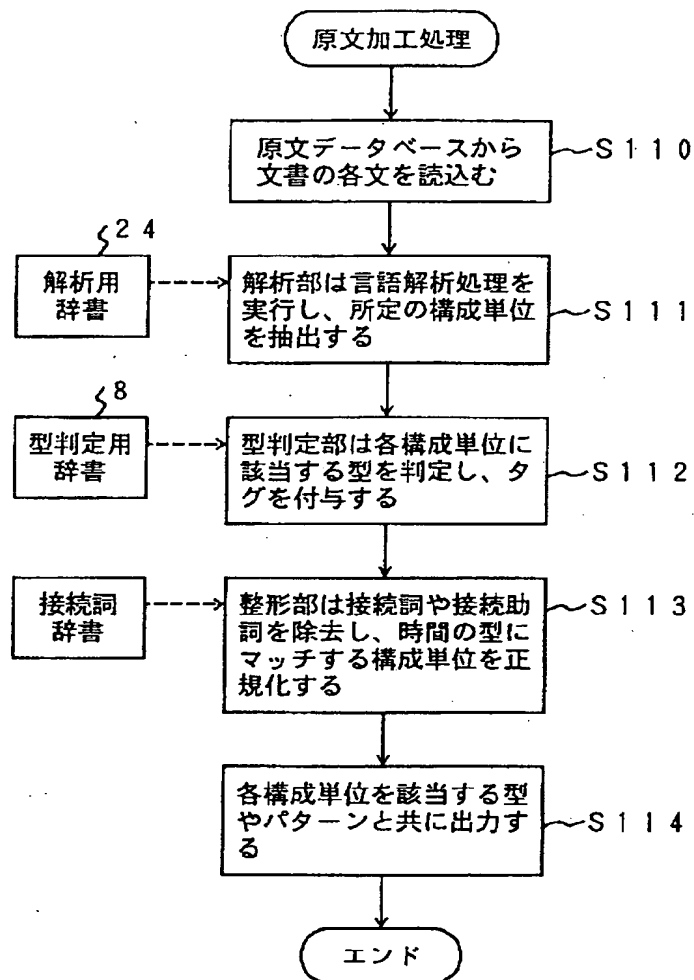
【 図29 】



【 図30 】



【 図3 1 】



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☐ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☒ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☒ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.